

Yash Akhauri



GitHub | LinkedIn | akhauri.yash@gmail.com | +91 78915 12802

A Research Scientist specializing in Neural Network optimization with a keen interest in Mathematics and Physics.

EDUCATION

BITS PILANI | B.E. IN ELECTRONICS AND INSTRUMENTATION

Aug 2016 - May 2020 | RJ, India

Fluent: Python3, PyTorch, Mathematica

Familiar: C++, Java, CUDA, OpenMP, Tensorflow, Android Studio, LibGDX, Docker

EXPERIENCE

INTEL | RESEARCH SCIENTIST

MAY 2020 - PRESENT, BANGALORE, INDIA

- Research Scientist at the Cloud Systems Research (CSR) Lab in the Systems and Software Research (SSR) Group.

XILINX RESEARCH | VISITING SCHOLAR

AUG 2019 - MAY 2020, DUBLIN, IRELAND

- Developed a library for co-design of neural network topologies and reconfigurable hardware that maps to an efficient FPGA implementation without the need for a custom accelerator architecture or a scheduler.
- Targeted the Jet Substructure Classification task as part of CERN LMS L1 trigger experiments, used the library to deploy models with 10x lower latency than FPGA4HEP designs. Demonstrated quantization library Brevitas to CERN [GitHub].

URANIOM | RESEARCH INTERN

JAN 2019 - JUL 2019, FRANCE (REMOTE)

- Implementing semantic segmentation models for face transfer across GIFs, progressive GANs for realistic UV map generation and explored effective weight-sharing strategies for neural networks under a research collaboration.

WOLFRAM | UNDERGRADUATE RESEARCHER

JUNE 2018 - JULY 2018, MASSACHUSETTS

- Developed HadaNet MLPs in the Wolfram Language and worked on C OpenMP kernels for GEMM and Convolutions using the Hadamard Binarization algorithm. [OpenMP kernel] | [CUDA kernel] | [Whitepaper]

RESEARCH

PAPERS

LOGICNETS: CO-DESIGNED NEURAL NETWORKS AND CIRCUITS FOR EXTREME-THROUGHPUT APPLICATIONS

[IEEE ¹] YASH AKHAURI*, YAMAN UMUROGLU*, NICHOLAS J. FRASER, MICHAELA BLOTT *FPL '20* | SWEDEN | SEPT 2020

Paper accepted at **FCCM'20** as a poster presentation. Video presentation can be found here.

HIGH-THROUGHPUT DNN INFERENCE WITH LOGICNETS

[IEEE ¹] YAMAN UMUROGLU, YASH AKHAURI, NICHOLAS J. FRASER, MICHAELA BLOTT *FCCM'20* | AR, USA | MAY 2020

HADANETS: FLEXIBLE QUANTIZATION STRATEGIES FOR NEURAL NETWORKS

[IEEE] YASH AKHAURI *CVPR'19 Workshop* | CA, USA | JUN 2019

Paper accepted at **CVPR'19 UAVision workshop - Orals**.

Delivered a "Theatre Talk" and poster at the Intel Demo Booth at CVPR'19.

EXPOSING HARDWARE BUILDING BLOCKS TO MACHINE LEARNING FRAMEWORKS

[ARXIV -- BACHELOR'S THESIS] YASH AKHAURI

DEC 2019

¹To be published

¹To be published

TALKS

WOLFRAM TECHNOLOGY CONFERENCE

SPEAKER

CHAMPAIGN, IL | OCT. 2018

Delivered a talk on my research on Hadamard Neural Networks.

INTEL AI MEETUP

SPEAKER

DELHI, IN | SEPT. 2018

[PPTX] [Article] Spoke about my research on scaling AI using Intel technologies. This event was organized by Intel.

INTEL AI DEVCON

POSTER

SAN FRANCISCO, BANGALORE | MAY & AUG 2018

Presented posters on quantized GEMM kernels for Intel Xeon Phi

GRANTS

INTEL NERVANA EARLY INNOVATORS GRANT

\$5000

Received research grant to develop Binary Precision Neural Networks and Real time Artistic Style Transfer. The technical article can be found [here.] The code can be found [here.]

Intel AI Academy Success Story [Link] published by Intel for the research done as part of this grant in the field of Quantized Neural Networks.

INTEL CVPR TRAVEL GRANT

\$3000

Received a travel grant from Intel to present research at the Intel Demo Booth at CVPR'19.

WOLFRAM STUDENT AID

\$2400

Received aid to attend the Wolfram Summer School and develop Hadamard Binary Neural Networks.

KVPY SCHOLAR

Selected as a KVPY scholar by the Department of Science and Technology, Government of India.

INSPIRE SCHOLARSHIP

Selected for the INSPIRE Scholarship by the Department of Science and Technology (DST), Government of India.

PROJECTS

EXPLOITING HUFFMAN CODING AND WEIGHT-SHARING FOR MEMORY-EFFICIENT INFERENCE ON FPGAs. PYTORCH

Developed a neural network weight sharing strategy and proposed a methodology to leverage this strategy in the FINN architecture for Multi-Layer Offload and Data-Flow Architectures of Neural Network deployment on FPGAs.

WHITEPAPER - IMPROVING DISTRIBUTED MESH COMPUTING WITH HADAMARD BINARY NEURAL NETWORKS.

[Link]

xGEMM & xCONV

C++, CUDA, OPENMP

Coded efficient 3D convolutional and GEMM kernels for XNOR (bit quantized) networks using CUDA C programming and OpenMP. Optimized kernels are for Intel processors and Nvidia GPUs. Invited to present a poster at Intel AI DevCon and Intel AI Student Ambassador Summit, San Francisco. The codes can be found here: [OpenMP kernel] | [CUDA kernel].

REAL TIME ARTISTIC STYLE TRANSFER

PYTHON, TENSORFLOW, OPENCV

/

GRAVDASH

JAVA, LIBGDX, ANDROID STUDIO

Developed an android game using Java and libGDX as the framework.

BLOG

PYTHON, JAVA, C++, TENSORFLOW, OPENMP

Maintaining a [blog] covering various AI related topics with over 10000 hits.